

# Original kmeans Slides

# Clustering: Grouping Related Docs



CS229: Machine Learning  
Carlos Guestrin  
Stanford University

Slides include content developed by and co-developed with  
Emily Fox

# Motivating clustering approaches

# Goal: Structure documents by topic

Discover groups (*clusters*) of related articles



SPORTS



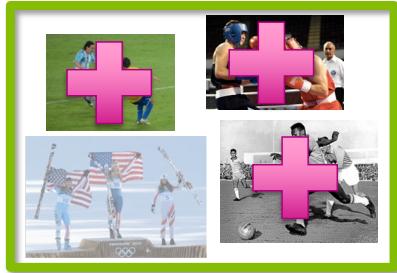
WORLD NEWS

# Why might clustering be useful?



# Learn user preferences

Set of clustered documents read by user



Cluster 1



Cluster 2



Cluster 3



Cluster 4



Use feedback  
to learn user  
preferences  
over topics

# Clustering: An unsupervised learning task

# What if some of the labels are known?

Training set of labeled docs



SPORTS



WORLD NEWS



ENTERTAINMENT



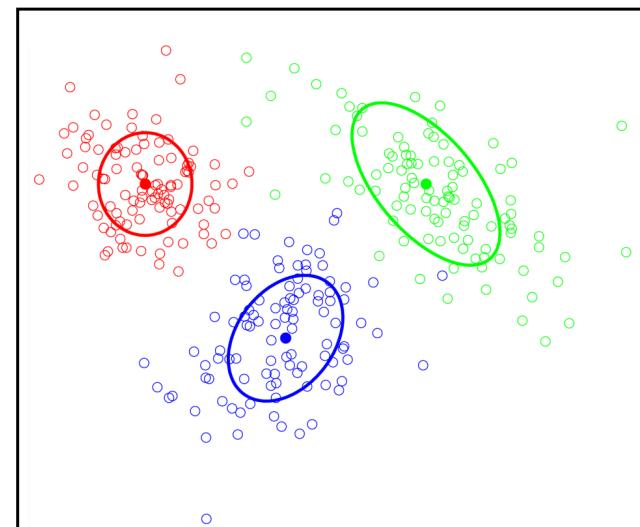
SCIENCE

# Clustering

No labels provided  
...uncover cluster structure  
from input alone

Input: docs as vectors  $x_i$   
Output: cluster labels  $z_i$

An unsupervised  
learning task

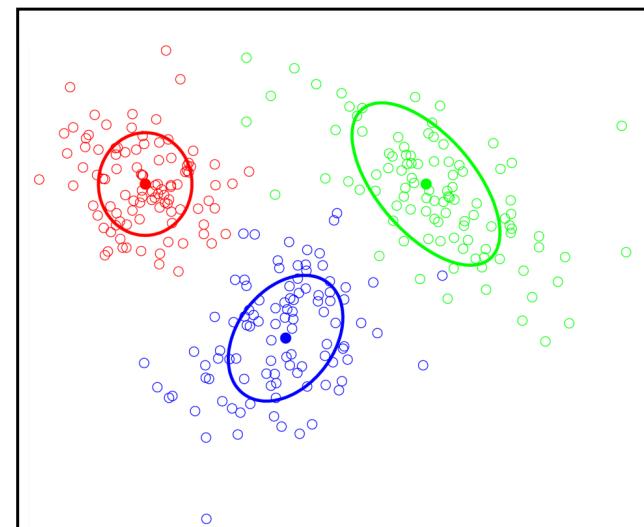


# What defines a cluster?

Cluster defined by **center** & **shape/spread**

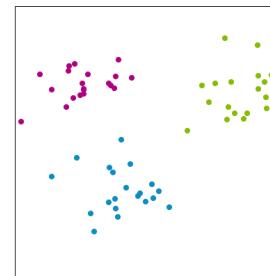
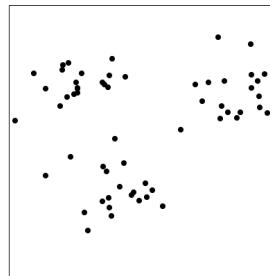
Assign observation  $x_i$  (**doc**)  
to cluster  $k$  (**topic label**) if

- Score under cluster  $k$  is higher than under others
- For simplicity, often define score as **distance to cluster center** (ignoring shape)

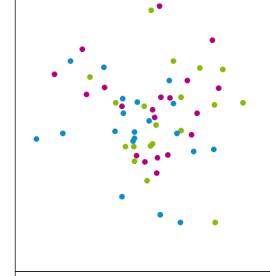
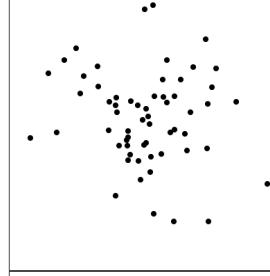


# Hope for unsupervised learning

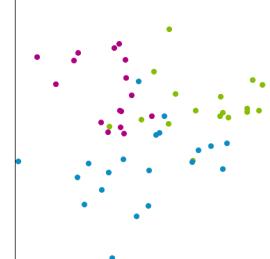
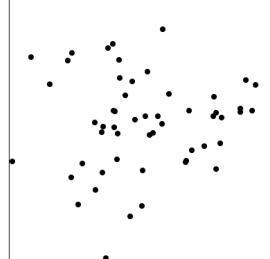
Easy



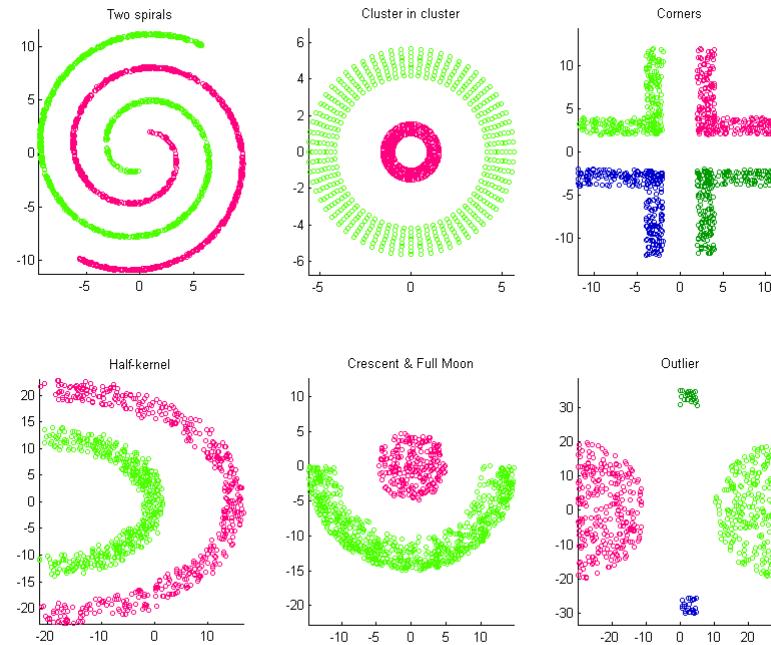
Impossible



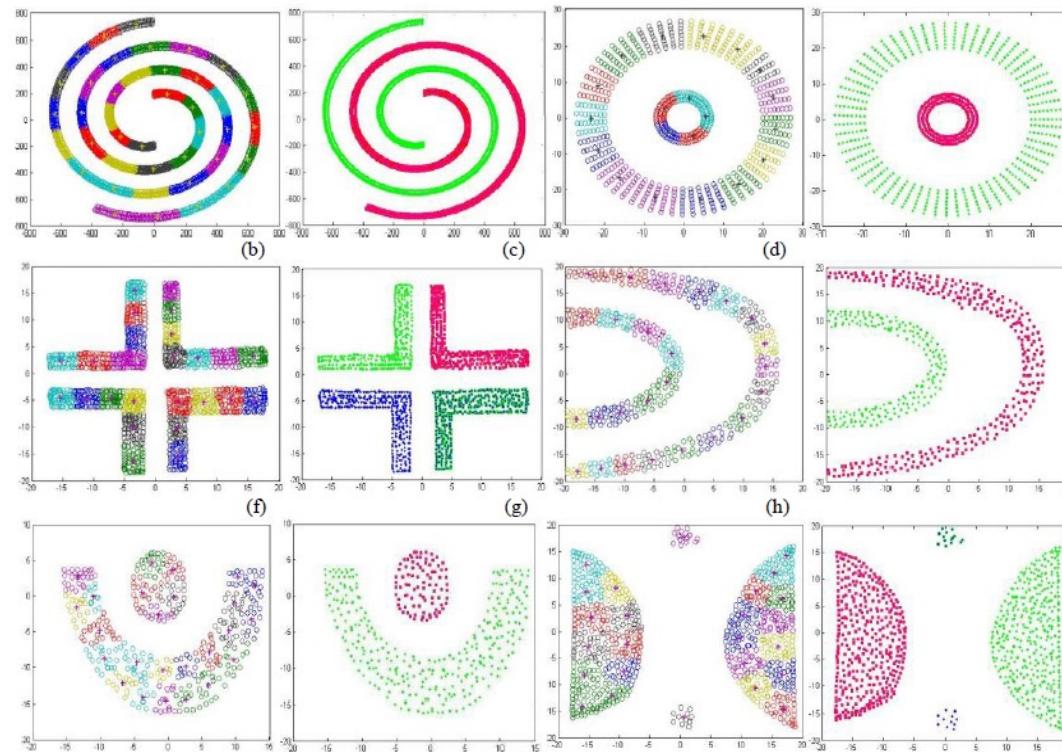
In between



# Other (challenging!) clusters to discover...



# Other (challenging!) clusters to discover...

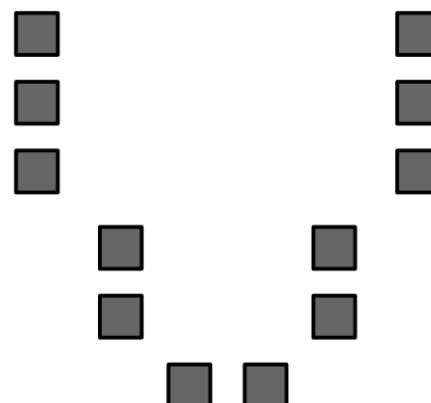


# k-means: A clustering algorithm

# k-means

Assume

- Score= distance to cluster center  
(smaller better)

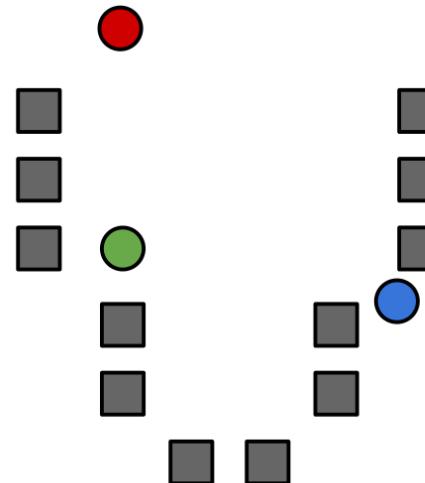


DATA  
to  
CLUSTER

# k-means algorithm

0. Initialize cluster centers

$$\mu_1, \mu_2, \dots, \mu_k$$

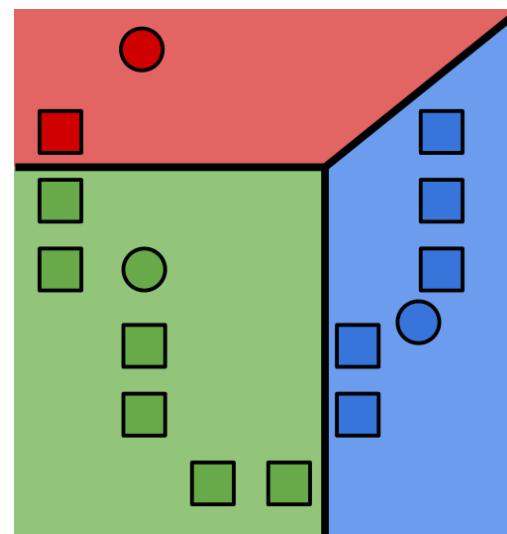


# k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

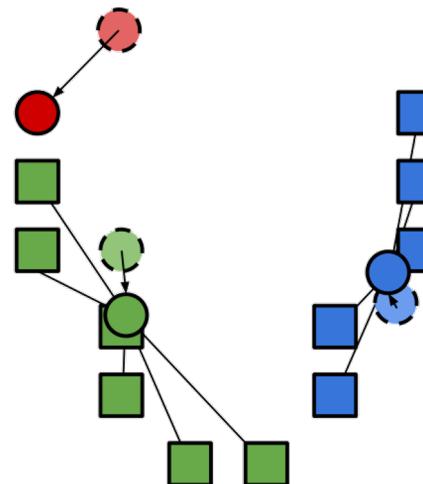
Inferred label for obs i, whereas  
supervised learning has given label  $y_i$



# k-means algorithm

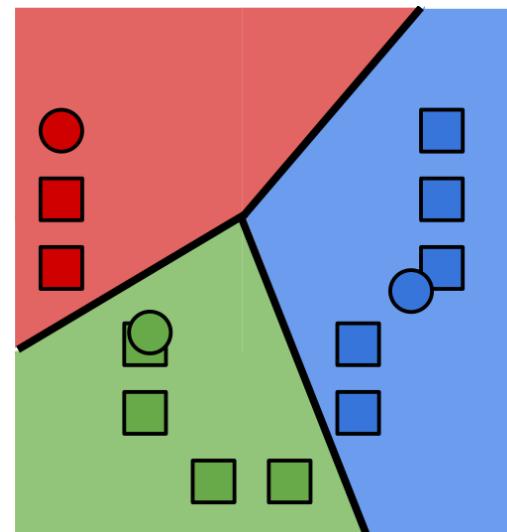
0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations

$$\mu_j = \frac{1}{n_j} \sum_{i:z_i=j} \mathbf{x}_i$$



# k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations
3. Repeat 1.+2. until convergence



# Why does K-means work???

- What's k-means optimizing?
- Does it always converge?

# What is k-means optimizing?

- Potential function  $F(\mu, z)$  of centers  $\mu$  and point allocations  $z$ :
- Optimal k-means:

# Does K-means converge??? Part 1

- Optimize potential function:

$$\min_{\mu} \min_{\mathbf{z}} F(\mu, \mathbf{z}) = \min_{\mu} \min_{\mathbf{z}} \sum_{j=1}^N \|\mu_{z_i} - x_i\|_2^2$$

- Fix  $\mu$  and minimize  $\mathbf{z}$ :

# Does K-means converge??? Part 2

- Optimize potential function:

$$\min_{\mu} \min_{\mathbf{z}} F(\mu, \mathbf{z}) = \min_{\mu} \min_{\mathbf{z}} \sum_{j=1}^N \|\mu_{z_i} - x_i\|_2^2$$

- Fix  $\mathbf{z}$  and minimize  $\mu$ :

# Coordinate descent algorithms

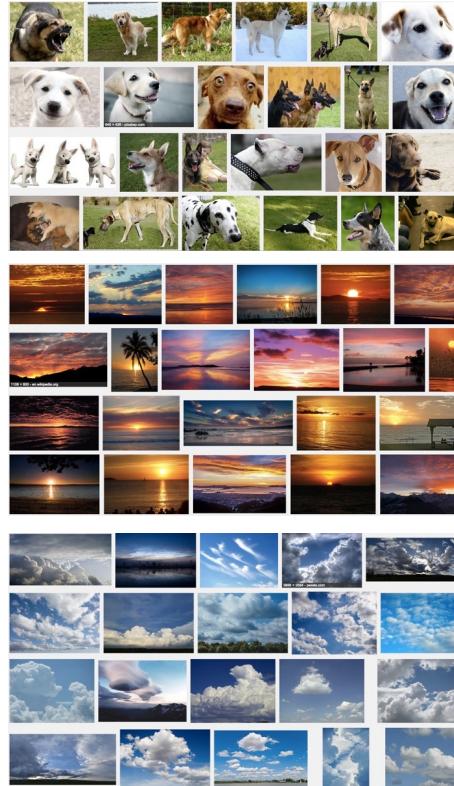
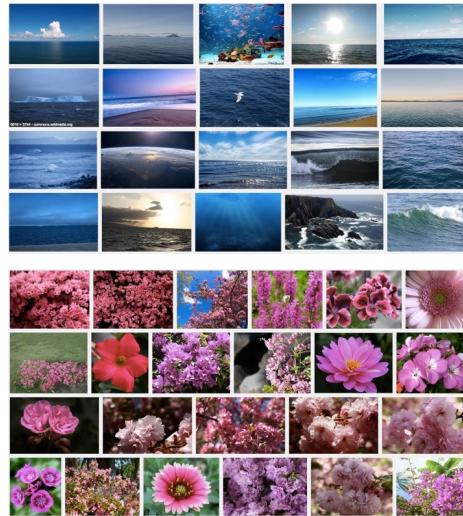
$$\min_{\mu} \min_{\mathbf{z}} F(\mu, \mathbf{z}) = \min_{\mu} \min_{\mathbf{z}} \sum_{j=1}^N \|\mu_{z_i} - x_i\|_2^2$$

- Want:  $\min_a \min_b F(a,b)$
- Coordinate descent:
  - fix a, minimize b
  - fix b, minimize a
  - repeat
- Converges!!!
  - if F is bounded
  - to a (often good) local optimum
    - as we saw in applet (play with it!)
      - (For LASSO it converged to the global optimum, because of convexity)
- K-means is a coordinate descent algorithm!

# Summary for k-means

# Clustering images

- For search, group as:
  - Ocean
  - Pink flower
  - Dog
  - Sunset
  - Clouds
  - ...



# Limitations of k-means

Assign observations to closest cluster center

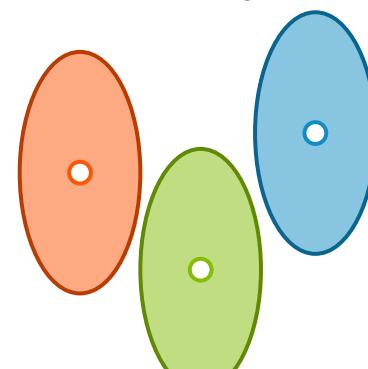
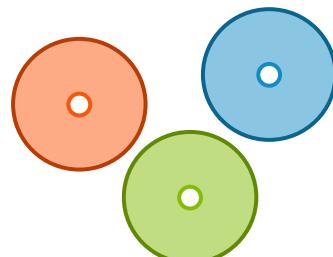
$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

Only center matters

Can use weighted Euclidean,  
but requires *known* weights

Still assumes all clusters have  
the same axis-aligned ellipses

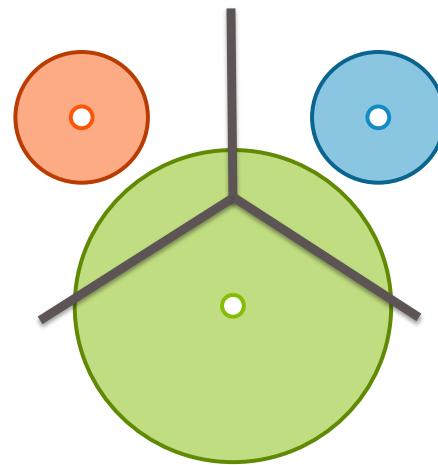
Equivalent to assuming  
*spherically symmetric* clusters



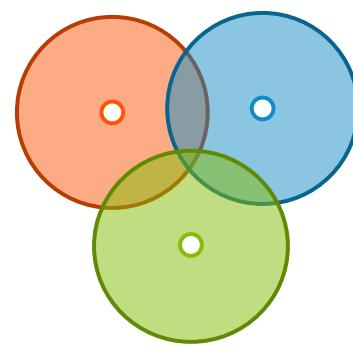
©2022 Carlos Guestrin

CS229: Machine Learning

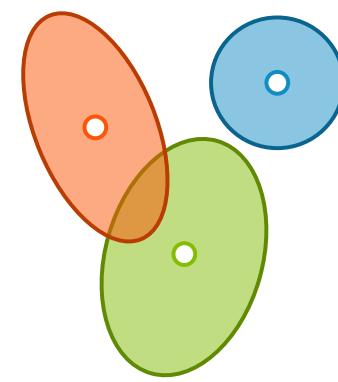
# Failure modes of k-means



disparate cluster sizes



overlapping clusters



different  
shaped/oriented  
clusters

# What you can do now...

- Describe the input (unlabeled observations) and output (labels) of a clustering algorithm
- Determine whether a task is supervised or unsupervised
- Cluster documents using k-means
- Describe potential applications of clustering

# Annotated kmeans Slides

# Clustering: Grouping Related Docs



CS229: Machine Learning  
Carlos Guestrin  
Stanford University

Slides include content developed by and co-developed with  
Emily Fox

# Motivating clustering approaches

# Goal: Structure documents by topic

Discover groups (*clusters*) of related articles



SPORTS



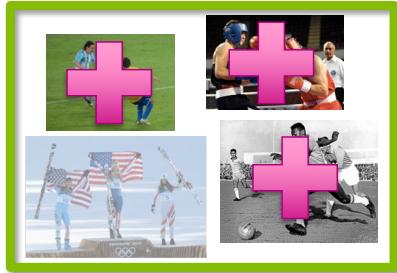
WORLD NEWS

# Why might clustering be useful?



# Learn user preferences

Set of clustered documents read by user



Cluster 1



Cluster 2



Cluster 3



Cluster 4



Use feedback  
to learn user  
preferences  
over topics

# Clustering: An unsupervised learning task

# What if some of the labels are known?

Training set of labeled docs



SPORTS



WORLD NEWS



ENTERTAINMENT



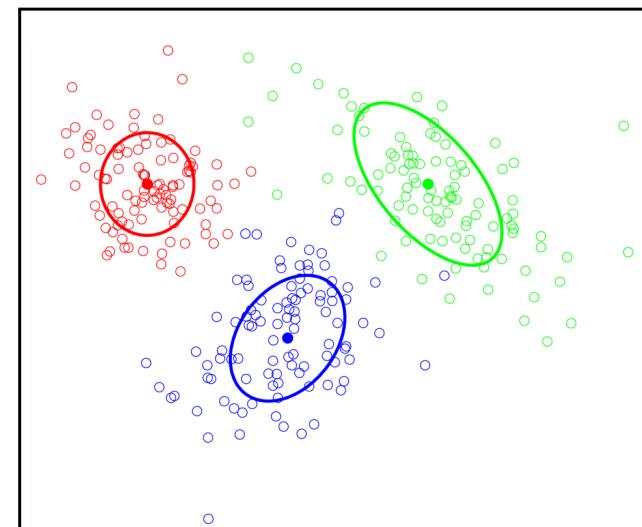
SCIENCE

# Clustering

No labels provided  
...uncover cluster structure  
from input alone

Input: docs as vectors  $x_i$   
Output: cluster labels  $z_i$

An unsupervised  
learning task

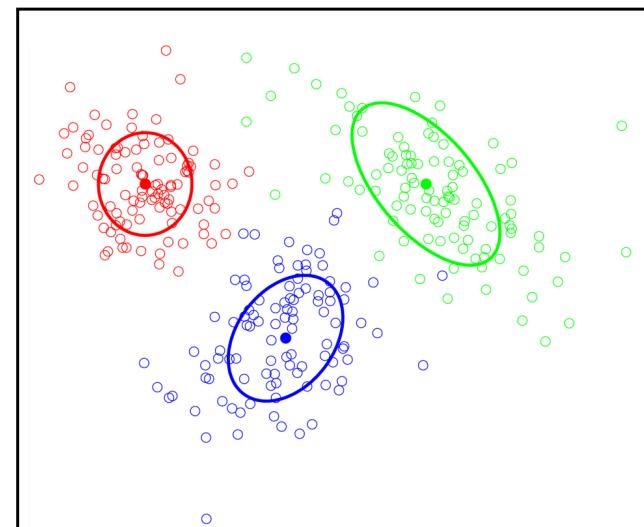


# What defines a cluster?

Cluster defined by **center** & **shape/spread**

Assign observation  $x_i$  (**doc**)  
to cluster  $k$  (**topic label**) if

- Score under cluster  $k$  is higher than under others
- For simplicity, often define score as **distance to cluster center** (ignoring shape)

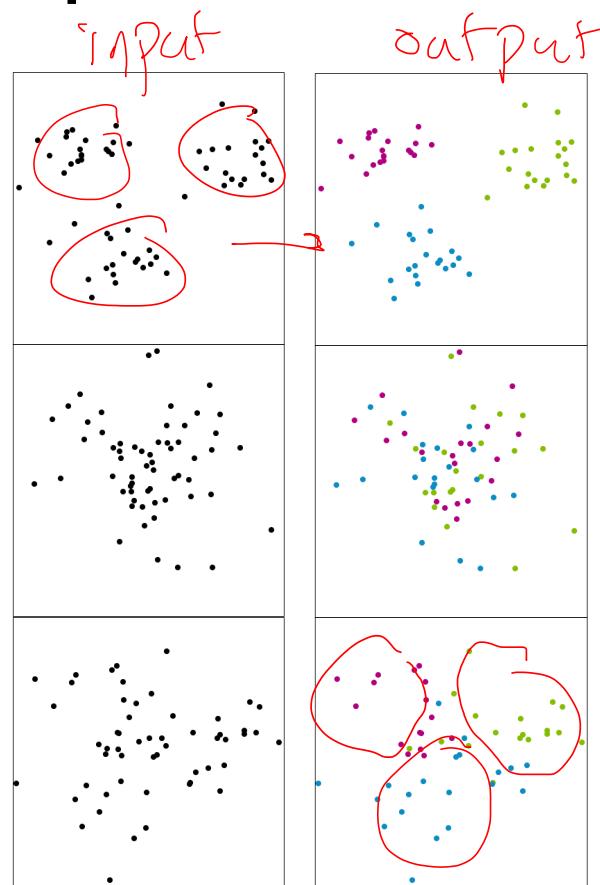


# Hope for unsupervised learning

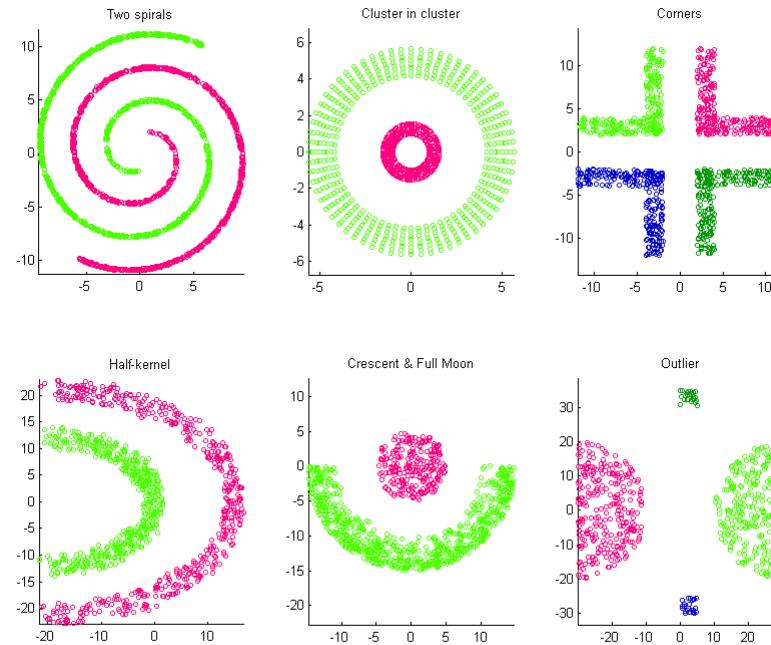
Easy

Impossible

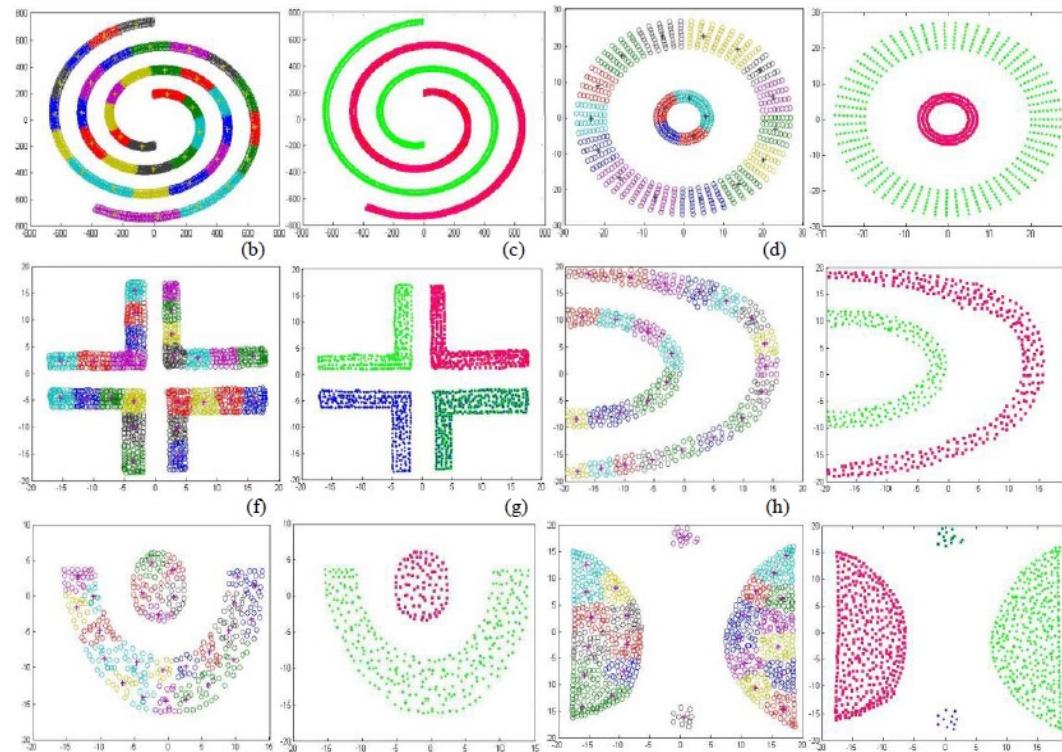
In between



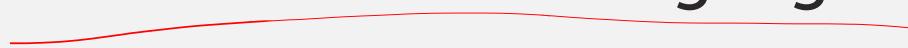
# Other (challenging!) clusters to discover...



# Other (challenging!) clusters to discover...



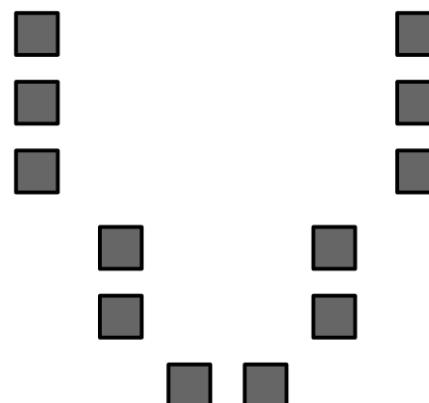
# k-means: A clustering algorithm



# k-means

Assume

- Score= distance to cluster center  
(smaller better)

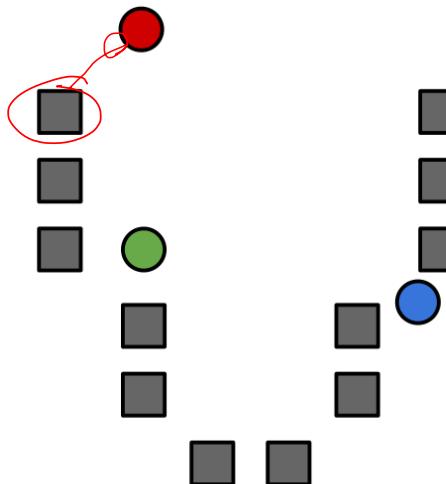


DATA  
to  
CLUSTER

# k-means algorithm

## 0. Initialize cluster centers

$$\mu_1, \mu_2, \dots, \mu_k$$



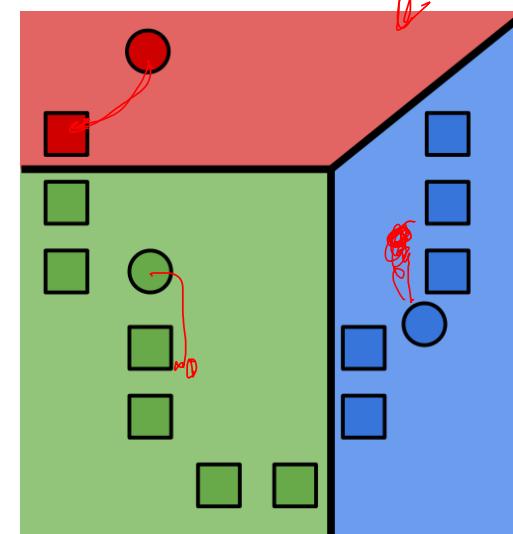
# k-means algorithm

0. Initialize cluster centers

1. Assign observations to  
closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

Inferred label for obs i, whereas  
supervised learning has given label  $y_i$



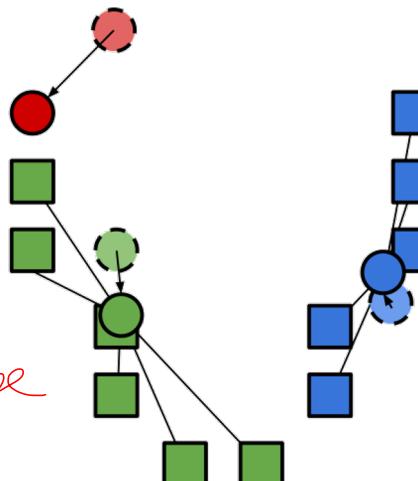
# k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations

boring center

$$\mu_j = \frac{1}{n_j} \sum_{i:z_i=j} \mathbf{x}_i$$

average

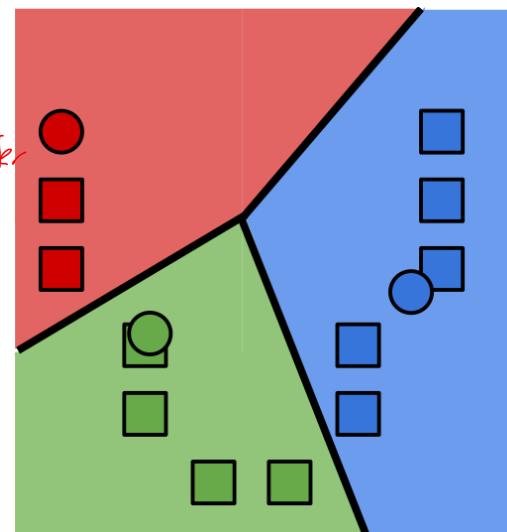


# k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations
3. Repeat 1.+2. until convergence

↳ no data point change  
Cluster assignment

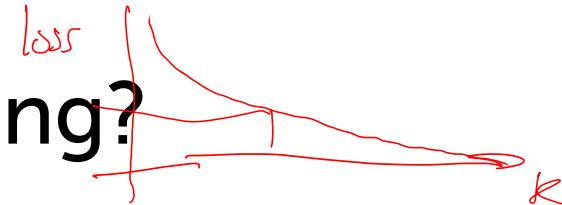
Classification Step



# Why does K-means work???

- What's k-means optimizing? → *What's the loss function*
- Does it always converge? ← *??*

# What is k-means optimizing?



- Potential function  $F(\mu, z)$  of centers  $\mu$  and point allocations  $z$ :

$$F(\mu, z) = \sum_{i=1}^N \|\mu_{z_i} - x_i\|_2^2$$

↗ loss function      ↗ cluster assignment      ↗ assigned cluster      ↗ input feature vector  
 $z_i \in \{1, \dots, k\}$

- Optimal k-means:

$$\min_{\mu} \min_z F(\mu, z)$$

# Does K-means converge??? Part 1

- Optimize potential function:

$$\min_{\mu} \min_{z} F(\mu, z) = \min_{\mu} \min_{z} \sum_{j=1}^N \|\mu_{z_i} - x_i\|_2^2$$

Classification step:

- Fix  $\mu$  and minimize  $z$ :

$$\min_{z_1, \dots, z_n} \sum_{i=1}^N \|\mu_{z_i} - x_i\|_2^2 = \sum_{i=1}^N \min_{\substack{z_i \in \{1, \dots, k\}}} \|\mu_{z_i} - x_i\|_2^2$$

Independent  $\min$  problem per data point  
≡ "classification step" in K-means

$\underbrace{\min_{z_1, \dots, z_n}}$   
no interaction terms between  $z_i$  &  $z_j$

# Does K-means converge??? Part 2

$n_j$  is Number of  
data points in cluster  
 $j$

- Optimize potential function:

$$\min_{\mu} \min_{z} F(\mu, z) = \min_{\mu} \min_{z} \sum_{j=1}^N \|\mu_{z_i} - x_i\|_2^2$$

*Recenter step of k-means*

*No interactions between  $\mu_i$ ,  $\mu_j$*

- Fix  $z$  and minimize  $\mu$ :

$$\begin{aligned} & \min_{\mu_1, \dots, \mu_K} \sum_{i=1}^N \|\mu_{z_i} - x_i\|_2^2 = \min_{\mu_1} \min_{\mu_2} \dots \min_{\mu_K} \sum_{j=1}^K \sum_{i: z_i=j} \|\mu_j - x_i\|_2^2 \\ &= \sum_{j=1}^K \underbrace{\min_{\mu_j} \sum_{i: z_i=j} \|\mu_j - x_i\|_2^2}_{\text{take derivative}} \end{aligned}$$

*Set to zero  $\Rightarrow$  optimal  $\mu_j$  is average:*

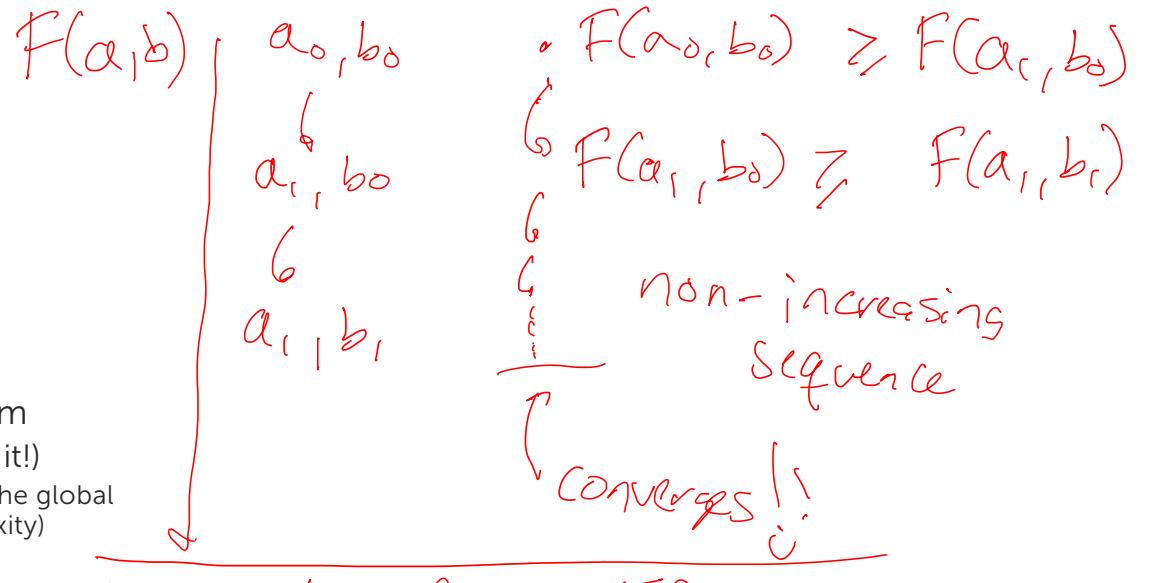
*find  $\mu_j$  that is closest on avg to points in cluster j*

$$\mu_j = \frac{\sum_{i: z_i=j} x_i}{n_j}$$

# Coordinate descent algorithms

$$\min_{\mu} \min_{\mathbf{z}} F(\mu, \mathbf{z}) = \min_{\mu} \min_{\mathbf{z}} \sum_{j=1}^N \|\mu z_i - x_i\|_2^2$$

- Want:  $\min_a \min_b F(a, b)$
- Coordinate descent:
  - fix  $a$ , minimize  $b$
  - fix  $b$ , minimize  $a$
  - repeat
- Converges!!!
  - if  $F$  is bounded
  - to a (often good) local optimum
    - as we saw in applet (play with it!)
    - (For LASSO it converged to the global optimum, because of convexity)
- K-means is a coordinate descent algorithm!

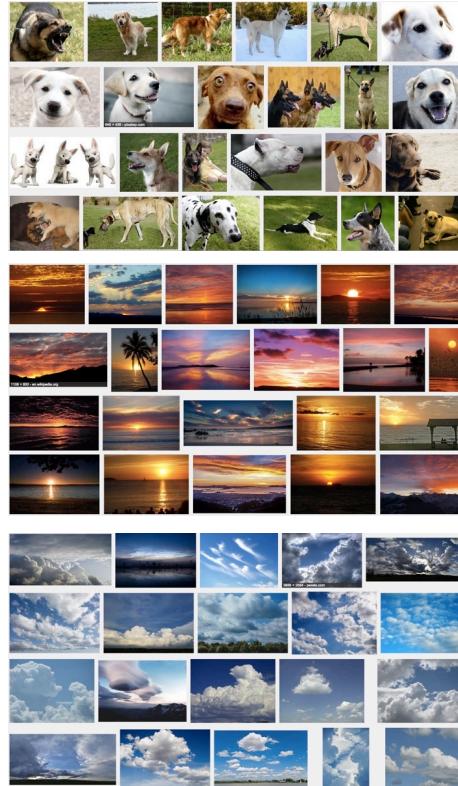
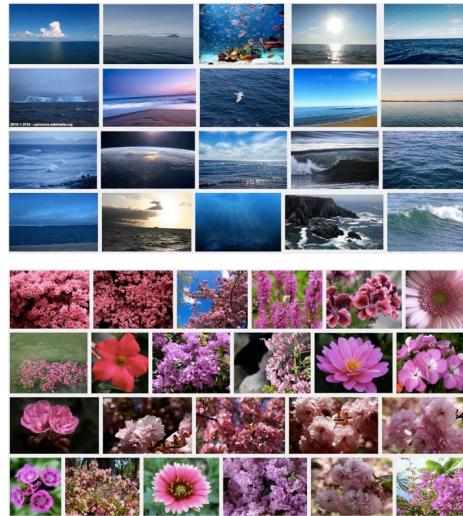


use random restarts lower bound  $\delta \nabla F(a, b)$   
or other tricks (k-means ++)

# Summary for k-means

# Clustering images

- For search, group as:
  - Ocean
  - Pink flower
  - Dog
  - Sunset
  - Clouds
  - ...



# Limitations of k-means

Assign observations to closest cluster center

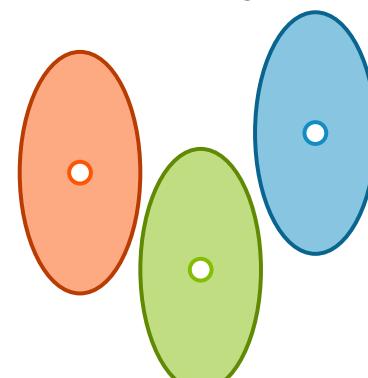
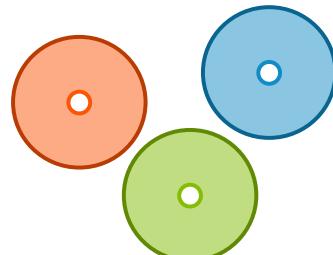
$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

Only center matters

Can use weighted Euclidean,  
but requires *known* weights

Still assumes all clusters have  
the same axis-aligned ellipses

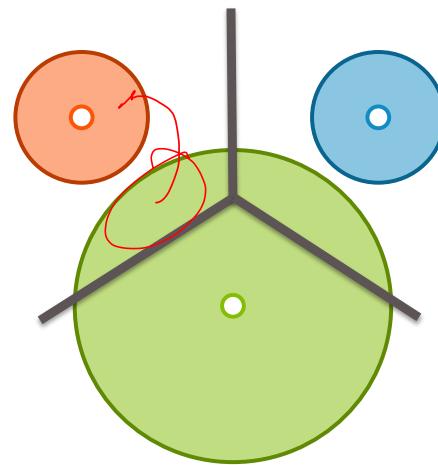
Equivalent to assuming  
*spherically symmetric* clusters



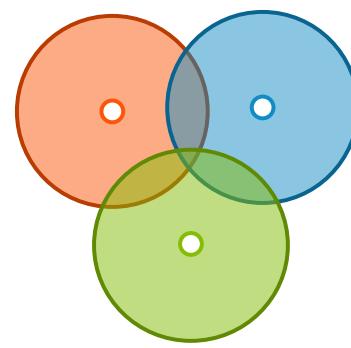
©2022 Carlos Guestrin

CS229: Machine Learning

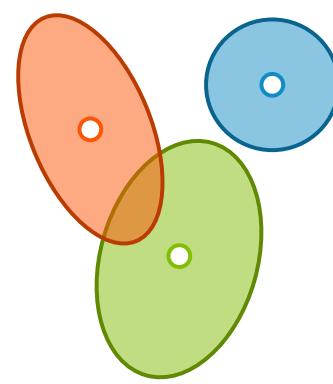
# Failure modes of k-means



disparate cluster sizes



overlapping clusters



different  
shaped/oriented  
clusters

# What you can do now...

- Describe the input (unlabeled observations) and output (labels) of a clustering algorithm
- Determine whether a task is supervised or unsupervised
- Cluster documents using k-means
- Describe potential applications of clustering